

HW# 6

Jiale Shi

1 Support vector machine

14.5

- A. The problem with kernelized ridge regression is that the solution vector w depends on all the training inputs. Vapnik (Vapnik et al. 1997) proposed a variant of the Huber loss function called the epsilon insensitive loss function, defined by

$$L_\epsilon \triangleq \begin{cases} 0 & \text{if } |y - \hat{y}| < \epsilon \\ |y - \hat{y}| - \epsilon & \text{otherwise} \end{cases} \quad (1)$$

The corresponding objective function is usually written in the following form

$$J = C \sum_{i=1}^N L_\epsilon(y_i, \hat{y}_i) + \frac{1}{2} \|w\|^2, \quad (2)$$

where

$$\hat{y}_i = w_0 + w^T x_i, \quad (3)$$

and $C = \frac{1}{\lambda}$ is a regularization constant. The objective is convex and unconstrained, but not differentiable, because of the absolute value function in the loss term. One popular approach is to formulate the problem as a constrained optimization problem. In particular, we introduce slack variables to represent the degree to which each point lies outside the tube:

$$y_i \leq f(x_i) + \xi_i^+ \quad (4)$$

$$y_i \geq f(x_i) - \xi_i^- \quad (5)$$

Given this, we can rewrite the objective as follows

$$J = C \sum_{i=1}^N (\xi_i^+ + \xi_i^-) + \frac{1}{2} \|w\|^2, \quad (6)$$

with this setup, the optimal solution has the form

$$\hat{w} = \sum_i \alpha_i x_i \quad (7)$$

where

$$\alpha_n = (a_n - \hat{a}_n) \quad (8)$$

where a_n and \hat{a}_n are the Lagrange multipliers corresponding to Eqs. 4 and 5.

- (a) For the regression support vector machine considered above, show that all training data points for which $\xi_n > 0$ will have $a_n = C$ and similarly all points for which $\hat{\xi}_n > 0$ will have $\hat{a}_n = C$.
 - (b) Compute the dual lagrangian for the support vector regression (Refer Eq. 7.61 of Bishop)
- B. For the data-set given [at this link](#), train a support vector machine with polynomial kernel p . Perform for various polynomial orders and plot order vs. error. To ensure hard margin, use $C = 10^6$.

(a). (6) must be minimized subject to the constraints

$$\xi_i \geq 0 \text{ and } \hat{\xi}_i \geq 0.$$

This can be achieved by introducing Lagrange multipliers

$$a_n \geq 0, \hat{a}_n \geq 0, \mu_n \geq 0, \hat{\mu}_n \geq 0$$

$$L = C \sum_{i=h}^N (\xi_i + \hat{\xi}_i) + \frac{1}{2} \|W\|^2 - \sum_{n=1}^N (\alpha_n \xi_n + \hat{\alpha}_n \hat{\xi}_n)$$

$$- \sum_{i=h}^N a_i (\epsilon + \xi_i + y_i - t_i) - \sum_{n=1}^N \hat{a}_n (\epsilon + \hat{\xi}_n - \hat{y}_n + t_n)$$

$$\frac{\partial L}{\partial W} = 0 \Rightarrow W = \underbrace{\sum_{n=1}^{\infty} (a_n - \hat{a}_n) \phi(x_n)}_{g = w^T \phi(x) + b}$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{n=1}^{\infty} (a_n - \hat{a}_n) = 0$$

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow \alpha_n + \mu_n = C \xrightarrow{\mu_n \geq 0} 0 \leq \alpha_n \leq C$$

$$\frac{\partial L}{\partial \hat{\xi}_n} = 0 \Rightarrow \hat{\alpha}_n + \hat{\mu}_n = C \xrightarrow{\hat{\mu}_n \geq 0} 0 \leq \hat{\alpha}_n \leq C$$

$$\tilde{L}(a, \hat{a}) = -\frac{1}{2} \sum_{h=1}^N \sum_{m=1}^N (a_h - \hat{a}_h)(a_m - \hat{a}_m) K(x_h, x_m)$$

$$- \epsilon \sum_{h=1}^N (a_h + \hat{a}_h) + \sum_{h=1}^N (a_h - \hat{a}_h) t_h$$

$$y(x) = \sum_{h=1}^N (a_h - \hat{a}_h) K(x, x_h) + b$$

$$a_h (\epsilon + \xi_h + y_h - t_h) = 0$$

$$\hat{a}_h (\epsilon + \hat{\xi}_h - \hat{y}_h + t_h) = 0$$

$$(C - a_h) \xi_h = 0, \quad \xi_h > 0, \quad a_h = C$$

$$(C - \hat{a}_h) \hat{\xi}_h = 0, \quad \hat{\xi}_h > 0, \quad \hat{a}_h = C$$

$$(b) \quad \tilde{L}(a, \hat{a}_n) = -\frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) k(x_n, x_m)$$

$$-\epsilon \sum_{n=1}^N (a_n + \hat{a}_n) + \sum_{n=1}^N (a_n - \hat{a}_n) t_n$$

subject to the constraints

$$0 \leq a_n \leq C/N$$

$$0 \leq \hat{a}_n \leq C/N$$

$$\sum_{n=1}^N (a_n - \hat{a}_n) = 0$$

$$\sum_{n=1}^N (a_n + \hat{a}_n) \leq C$$

Proof:

$$\begin{aligned} L &= C \sum_{n=1}^N (\hat{\zeta}_n + \hat{\hat{\zeta}}_n) + \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (\mu_n \hat{\zeta}_n + \hat{\mu}_n \hat{\hat{\zeta}}_n) \\ &\quad - \sum_{n=1}^N a_n (\epsilon + \hat{\zeta}_n + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (\epsilon + \hat{\hat{\zeta}}_n - y_n + t_n) \\ &= C \sum_{n=1}^N (\hat{\zeta}_n + \hat{\hat{\zeta}}_n) + \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (a_n + \mu_n) \hat{\zeta}_n - \sum_{n=1}^N (\hat{a}_n + \hat{\mu}_n) \hat{\hat{\zeta}}_n \\ &\quad - \sum_{n=1}^N a_n (\epsilon + y_n - t_n) - \sum_{n=1}^N \hat{a}_n (\epsilon - y_n + t_n) \\ &= C \sum_{n=1}^N (\hat{\zeta}_n + \hat{\hat{\zeta}}_n) + \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (\hat{\zeta}_n - \sum_{n=1}^N (a_n + \mu_n) \hat{\zeta}_n) \\ &\quad - \sum_{n=1}^N (a_n + \hat{a}_n) \in - \sum_{n=1}^N (a_n - \hat{a}_n)(y_n - t_n) \end{aligned}$$

$$= \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (a_n + \hat{a}_n) \epsilon - \sum_{n=1}^N (a_n - \hat{a}_n)(y_n - t_n)$$

$$= \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (a_n - \hat{a}_n)(w^\top \phi(x_n) + b - t_n) - \sum_{n=1}^N (a_n + \hat{a}_n) \epsilon$$

$$= \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (a_n - \hat{a}_n)(w^\top \phi(x_n) + b) - \sum_{n=1}^N (a_n + \hat{a}_n) \epsilon + \sum_{n=1}^N (a_n - \hat{a}_n) t_n$$

$$= \frac{1}{2} \|w\|^2 - \sum_{n=1}^N (a_n - \hat{a}_n) w^\top \phi(x_n) - \sum_{n=1}^N (a_n + \hat{a}_n) \epsilon + \sum_{n=1}^N (a_n - \hat{a}_n) t_n$$

$$= \frac{1}{2} \|w\|^2 - \|w\|^2 - \sum_{n=1}^N (a_n + \hat{a}_n) \epsilon + \sum_{n=1}^N (a_n - \hat{a}_n) t_n$$

$$= -\frac{1}{2} \|w\|^2 - \sum_{n=1}^N (a_n + \hat{a}_n) \epsilon + \sum_{n=1}^N (a_n - \hat{a}_n) t_n$$

$$W = \sum_{n=1}^N (a_n - \hat{a}_n) \phi(x_n)$$

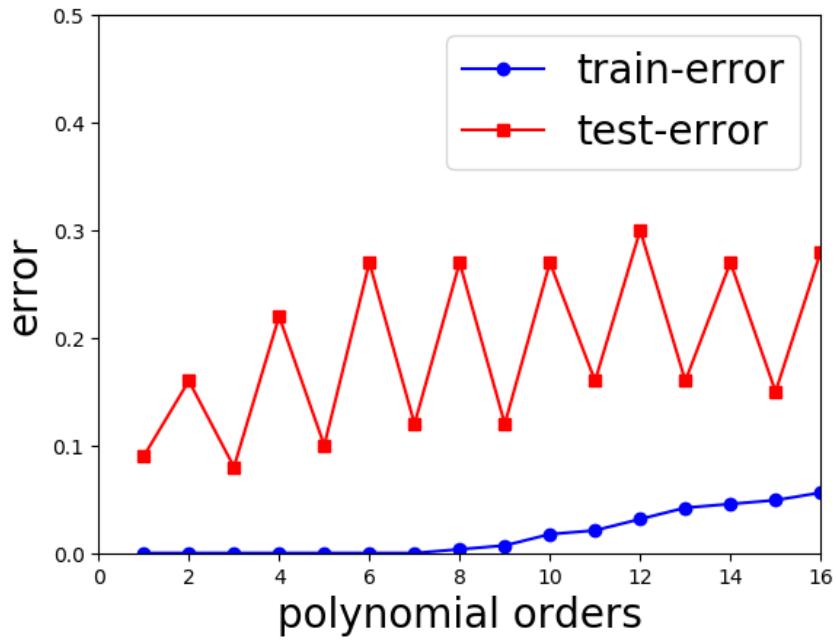
$$\tilde{\Sigma}(a, \hat{a}) = -\frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) \phi(x_n)^\top \phi(x_m)$$

$$+ \epsilon \sum_{n=1}^N (a_n + \hat{a}_n) - \sum_{n=1}^N (a_n - \hat{a}_n) t_n$$

$$= -\frac{1}{2} \sum_{m=1}^N \sum_{n=1}^N (a_n - \hat{a}_n)(a_m - \hat{a}_m) K(x_n, x_m)$$

$$+ \epsilon \sum_{n=1}^N (a_n + \hat{a}_n) - \sum_{n=1}^N (a_n - \hat{a}_n) t_n$$

B.



Error = 1 - accuracy.

use the `sklearn.SVM`

see code in the attached file

2 Relevance vector machine

14.3.2

- A. The relevance vector machine for regression is a linear model but with a modified prior that results in sparse solutions. The model defines a conditional distribution for a real-valued target variable t given an input vector \mathbf{x} , which takes the form

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}), \beta^{-1}) \quad (9)$$

The symbols in above equation takes the usual form.

Suppose we are given a set of N observations of the input vector \mathbf{x} , which we denote collectively by a data matrix \mathbf{X} whose n -th row is \mathbf{x}_n^T with $n = 1, \dots, N$. The corresponding target values are given by $\mathbf{t} = (t_1, \dots, t_N)^T$. Thus, the likelihood function is given by

$$p(t|\mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N p(t_i|\mathbf{x}_i, \mathbf{w}, \beta) \quad (10)$$

Next we introduce a prior distribution over the parameter vector \mathbf{w} as

$$p(\mathbf{w}|\alpha) = \prod_{i=1}^M \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha_i^{-1}), \quad (11)$$

where α_i represents the precision of the corresponding parameter w_i . When we maximize the evidence with respect to these hyperparameters, a significant proportion of them go to infinity, and the corresponding weight parameters have posterior distributions that are concentrated at zero.

- (a) For RVM discussed above, compute mean and covariance of the posterior distribution over weights.
- (b) Derive the result for the marginal likelihood function in the regression RVM, by performing the Gaussian integral over \mathbf{w} in (7.84) using the technique of completing the square in the exponential.
- (c) Derive the re-estimation equations as discussed in the notes.

- B. For the data-set provided [at this link](#), train a regression model by using RVM. Use

$$y(\mathbf{x}) = \sum_{i=1}^N w_i k(\mathbf{x}, \mathbf{x}_i) + b \quad (12)$$

where b is the bias parameter. Consider the kernel k to be Gaussian with kernel width 5.5. Plot the predictions along with proper prediction bounds.

A.
(a)

$$P(w|\alpha, \beta, \mathbf{z}) \propto P(w|\alpha) P(t|\mathbf{x}, w, \beta)$$

$$P(w|\alpha) = \prod_{j=1}^M P(w_j|\alpha_j) = \prod_{j=1}^M N(w_j|0, \alpha_j^{-1}) = \mathcal{N}_M(w|\text{diag}(\alpha)^{-1})$$

$$P(t|\mathbf{x}, w, \beta) = \prod_{i=1}^N N(t^{(i)} | w^T \phi(x^{(i)}), \beta^{-1}) = \mathcal{N}_N(t|\Phi w, \beta^{-1} I_N)$$

$$\sum_{i=1}^N \|w^T \phi(x^{(i)}) - t^{(i)}\|^2 = (\Phi w - \mathbf{y})^T (\Phi w - \mathbf{y})$$

$$(\bar{\Phi}w - t)^T (\beta I_N) (\bar{\Phi}w - t) + W^T \text{diag}(\lambda) W$$

$$= W^T (\beta \bar{\Phi}^T \bar{\Phi}) W - 2W^T (\beta \bar{\Phi}^T) t + \beta t^T t + W^T \text{diag}(\lambda) W$$

$$= W^T (\text{diag}(\lambda) + \beta \bar{\Phi}^T \bar{\Phi}) W - 2W^T (\beta \bar{\Phi}^T) t$$

$$= W^T \Sigma^{-1} W - 2W^T \Sigma^{-1} (\beta \Sigma \bar{\Phi}^T) t$$

$$\Sigma = (\text{diag}(\lambda) + \beta \bar{\Phi}^T \bar{\Phi})^{-1}$$

$$m = \beta \Sigma \bar{\Phi}^T t$$

(b) Marginal likelihood function obtained by integrating out the weight parameters

$$P(t|X, \lambda, \beta) = \int P(t|X, w, \beta) P(w|\lambda) dW$$

Because this represents the convolution of two Gaussians, it is readily evaluated to give the log marginal likelihood

$$\log P(t|X, \lambda, \beta) = \ln N(t|\mu, C)$$

$$C \in R^{N \times N}$$

$$C = \beta^{-1} I_N + \Phi \text{diag}(\lambda)^{-1} \Phi^T$$

proof: Using Results for linear Gaussian Models

$$P(x) = N(x | \mu, \Lambda^{-1})$$

$$P(y|x) = N(y | Ax + b, L^{-1})$$

For the above linear Gaussian model, the following hold.

$$P(y) = N(y | A\mu + b, L^{-1} + AA^{-1}A^T)$$

for this problem $x \rightarrow w, \mu \rightarrow 0, \Lambda \rightarrow \text{diag}(\lambda), y \rightarrow t, A \rightarrow \Phi$

$$b = 0, L = \beta I_N$$

$$P(t | \lambda, \beta) = N(t | 0, \beta^{-1} I_N + \Phi \text{diag}(\lambda)^{-1} \Phi^T)$$

(C)

$$P(D|\alpha, \beta) = (2\pi)^{-\frac{N}{2}} \beta^{\frac{N}{2}} |\text{diag}(\alpha)|^{\frac{1}{2}} \exp\left[-\frac{1}{2} t^T (\beta I_N - \beta^2 \Phi \Sigma \Phi^T)^{-1} t\right] |\bar{z}|^{\frac{1}{2}}$$

where $\Sigma = (\text{diag}(\alpha) + \beta \Phi^T \Phi)^{-1}$, $m = \beta \bar{z} \Phi^T t$

$$t^T \beta^2 \Phi \Sigma \Phi^T t = m^T \Sigma^{-1} m$$

$$P(D|\alpha, \beta) = (2\pi)^{-\frac{N}{2}} \beta^{\frac{N}{2}} |\text{diag}(\alpha)|^{\frac{1}{2}} \exp\left[-\frac{1}{2} (\beta t^T t - m^T \Sigma^{-1} m)\right] |\bar{z}|^{\frac{1}{2}}$$

$$\begin{aligned} \ln P(D|\alpha, \beta) &= -\frac{N}{2} \ln(2\pi) + \frac{N}{2} \ln \beta + \frac{1}{2} \ln |\text{diag}(\alpha)| \\ &\quad - \frac{1}{2} (\beta t^T t - m^T \Sigma^{-1} m) + \frac{1}{2} \ln |\bar{z}| \end{aligned}$$

$$\frac{\partial \ln P(D|\alpha, \beta)}{\partial \alpha_i} = 0$$

$$\frac{\partial \ln |\Sigma|}{\partial \alpha_i} = \text{Tr}\left(\Sigma + \frac{\partial \Sigma}{\partial \alpha_i}\right) = -\text{Tr}\left(\Sigma \frac{\partial \Sigma^{-1}}{\partial \alpha_i}\right) = -\bar{\Sigma}_{ii}$$

$$m^T \Sigma^{-1} \frac{\partial m}{\partial \alpha_i} = -m^T \Sigma^+ \Sigma \beta \frac{\partial \bar{\Sigma}^{-1}}{\partial \alpha_i} \Phi^T t = -m_i^2$$

$$\frac{1}{2\alpha_i} + \frac{1}{2} m_i^2 - m_i^2 - \frac{1}{2} \bar{\Sigma}_{ii} = 0$$

$$\boxed{\alpha_i = \frac{1 - \alpha_i \bar{\Sigma}_{ii}}{m_i^2} \equiv \frac{\gamma_i}{m_i^2}}$$

$$\frac{\partial \ln P(D|\alpha, \beta)}{\partial \beta} = 0$$

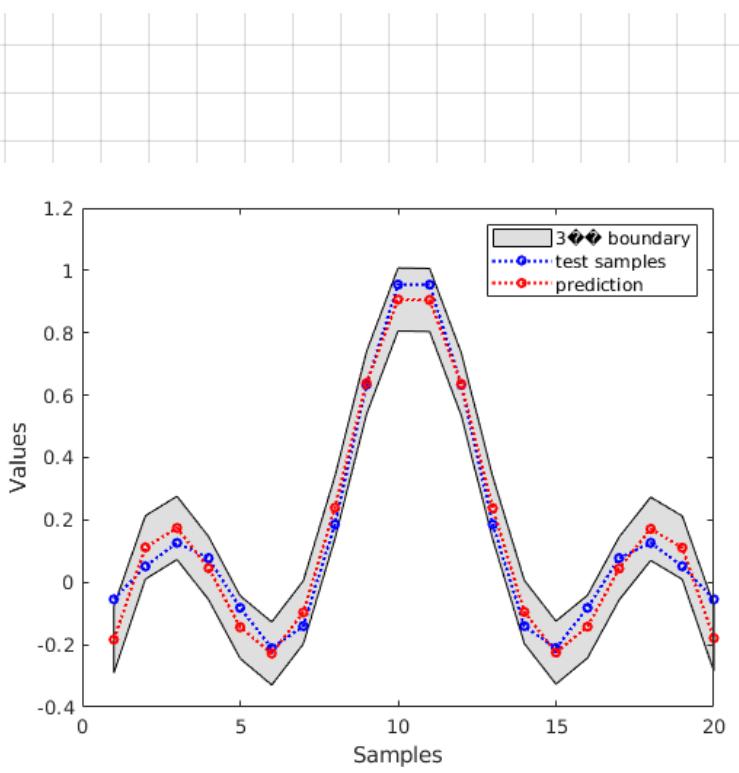
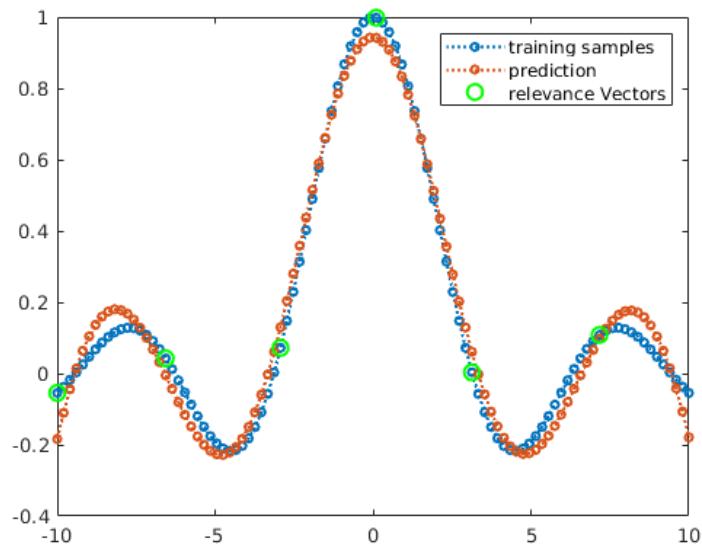
$$\frac{\partial \ln |\Sigma|}{\partial \beta} = \text{Tr} \left(\Sigma^{-1} \frac{\partial \Sigma}{\partial \beta} \right) = -\text{Tr} \left(\Sigma \frac{\partial \Sigma^{-1}}{\partial \beta} \right) = -\text{Tr} (\Sigma \Phi^T \Phi)$$

$$m = \beta \Sigma \Phi^T t$$

$$\frac{N}{2\beta} - \frac{1}{2} t^T t + \frac{1}{2} m^T \Phi^T \Phi m + m^T \Phi^T t - \frac{1}{2} \text{tr} (\Sigma \Phi^T \Phi) = 0$$

$$\boxed{(\beta^{(\text{new})})^{-1} = \frac{\|t - \Phi m\|^2}{N - \sum Y_i}}$$

B.



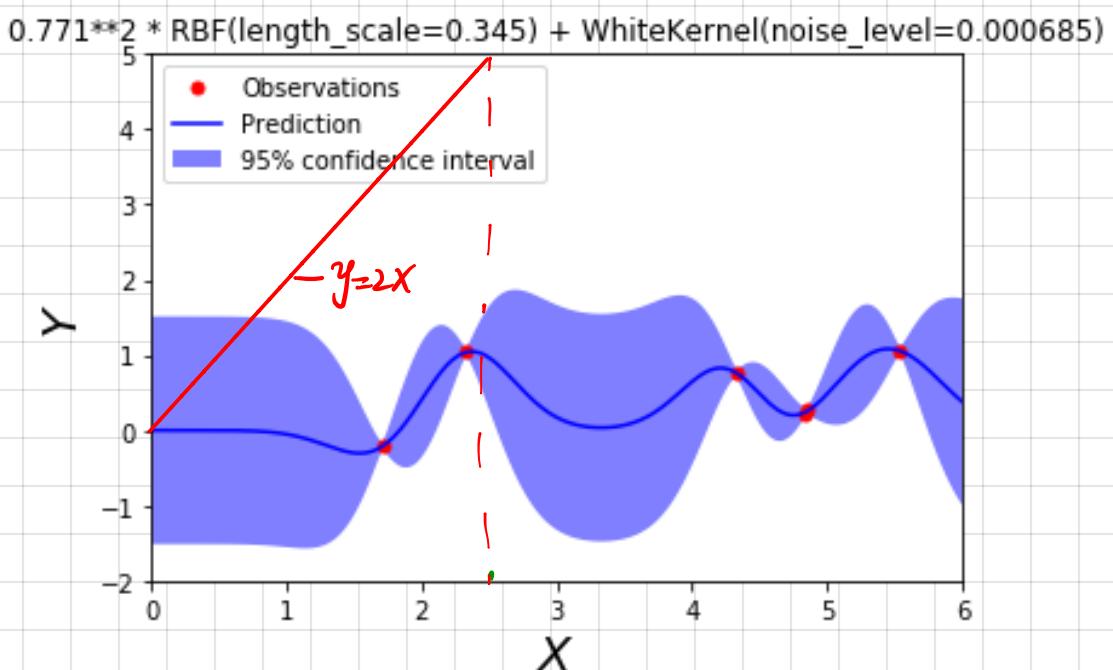
3 Gaussian process

Six training points are given, with the input $X = [1.71, 2.33, 4.33, 4.84, 4.86, 5.54]$ and the output $Y = [-0.2138, 1.0389, 0.7630, 0.2271, 0.2733, 1.0565]$. Implement the zero-mean Gaussian process regression model with the following covariance kernel:

$$k(x, x') = \theta_1 \exp\left(-\frac{(x - x')^2}{2\theta_2}\right) + \theta_3 \delta(x, x') \quad (13)$$

where δ is the kronecker delta.

- A. Learn the values of θ_{1-3} by maximizing the log-likelihood (MLE) based on the training points. The log-likelihood is known to be non-convex, so try several starting points for optimization to see if that affects the results.
- B. Run your Gaussian process with the learnt parameters to predict for the testing points on the interval $x^* \in [1, 6]$. Plot the prediction mean as well as the 95% confidence interval (two-sigma) of the learnt Gaussian process, and compare with their exact values computed by the formula: $y = \sim(2x)$. Comment on the prediction performance.



$$\theta_1 = 0.771^2$$

$$\theta_2 = (0.345)^2$$

$$\theta_3 = 0.000685$$

I try several starting points for optimization,

It won't affect the results.

I use python package `sklearn.gaussian_process`

The kernel's expression

$$\text{kernel} = \text{ConstantKernel} * RBF + \text{WhiteKernel}$$

B, see the plot.

The predict value is quit away from their exact formula: $y = \sim(2x)$

The prediction performance is bad.

The reason is the
the amount of input data is small,
the input data is quite away from the
exact formula : $y = \sim(2x)$.

4 Gaussian process latent variable model

- A. One approach to unsupervised learning with GPs is the Gaussian process latent variable model (GP-LVM) proposed by Lawrence (2004, 2005). GP-LVM can be considered as a multiple-output GP regression model where only the output data are given.

The inputs are unobserved and are treated as latent variables, however instead of integrating out the latent variables, they are optimized. This trick makes the model tractable and some theoretical grounding for the approach is given by the fact that the model can be seen as a nonlinear extension of the linear probabilistic PCA (PPCA).

Considering, $\mathbf{Y} \in \mathbb{R}^{N \times D}$ be the observed data and $\mathbf{X} \in \mathbb{R}^{N \times Q}$ to be the associated latent variables where $D > Q$,

- (a) provide an expression for the likelihood of the data. Assume, the GPs to be independent across the features.
- (b) Considering a prior of the form

$$p(\mathbf{X}) = \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n | \mathbf{0}, \mathbf{I}_Q), \quad (14)$$

The marginal likelihood of the data can be written as

$$p(\mathbf{Y}) = \int p(\mathbf{X}) p(\mathbf{Y}|\mathbf{X}) d\mathbf{X} \quad (15)$$

The marginal distribution shown in Eq. 15 is not tractable. Why?

- (c) Provide a variational approximation for the marginal distribution.

- B. For the oil data given [at this link](#), use GPLVM to identify the latent dimensions. Consider latent dimensions to be 3. Show the three latent dimensions.

(a)

$$P(Y_n | X_n, W, \beta) = N(Y_n | W X_n, \beta^{-1} I)$$

(b)

$$P(Y) = \int P(X) P(Y|X) dX$$

The marginal distribution in Eq. 15 is not tractable.

because only the output data Y is given. We don't know the input data X

(C) Instead, what we can do is

the dual approach of marginalising W and optimizing each X_n .

Integrating Parameters, Optimising Latent Variables.

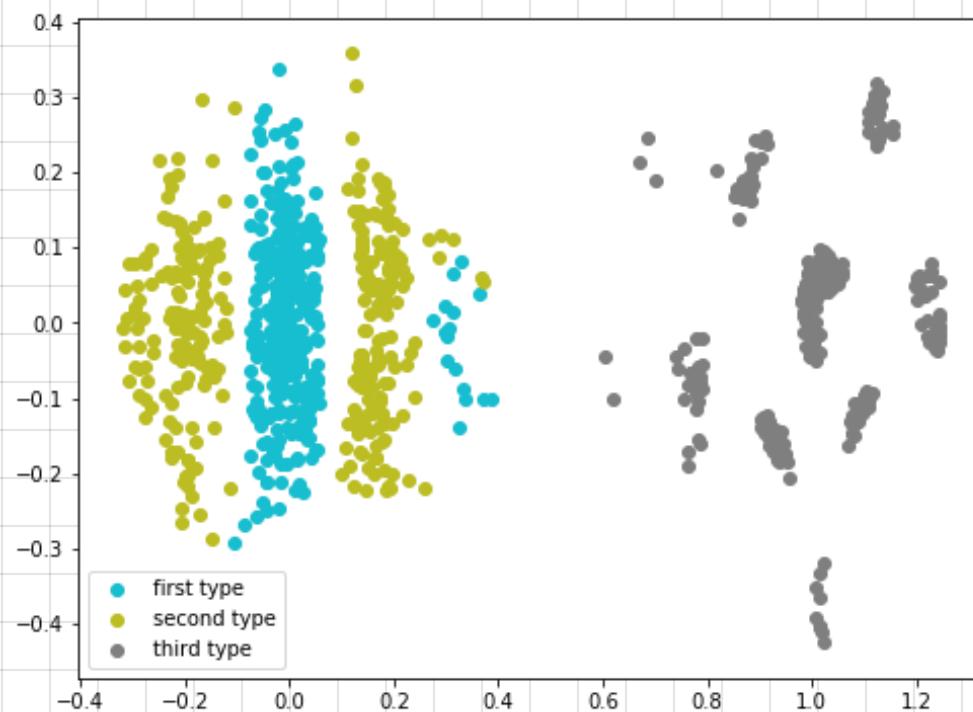
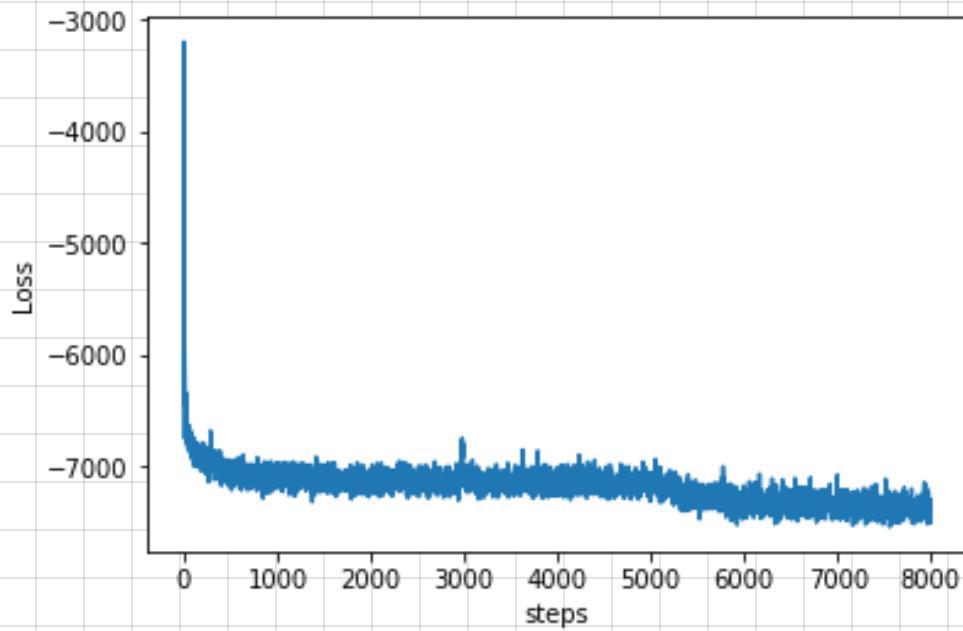
$$P(W) = \prod_{i=1}^D N(w_i | 0, \alpha^{-1} I)$$

$$P(Y|X, \beta) = \frac{1}{(2\pi)^{\frac{DN}{2}} |K|^{\frac{D}{2}}} \exp\left(-\frac{1}{2} \text{tr}(K^{-1} Y Y^T)\right)$$

$$\text{where } K = \alpha X X^T + \beta^{-1} I$$

B.

I use the train data and train label.



5 Boosting

An alternative approach to the kernel methods is to learn useful features directly from the input data. With this, we can create what we call an adaptive basis function model (ABM), which is a model of the form

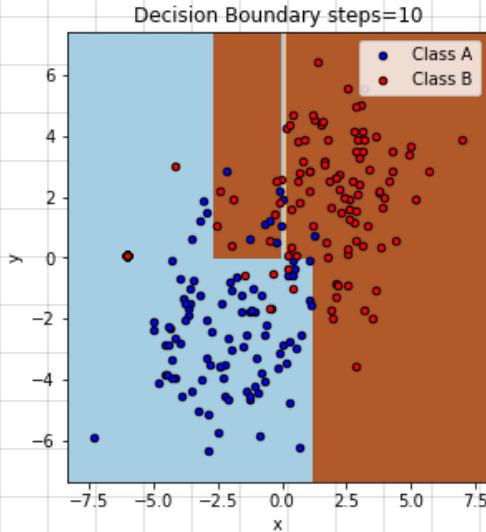
$$f(\mathbf{x}) = w_0 + \sum_{m=1}^M w_m \phi_m(\mathbf{x}), \quad (16)$$

where $\phi_m(\mathbf{x})$ is the mth basis function, which is learned from data.

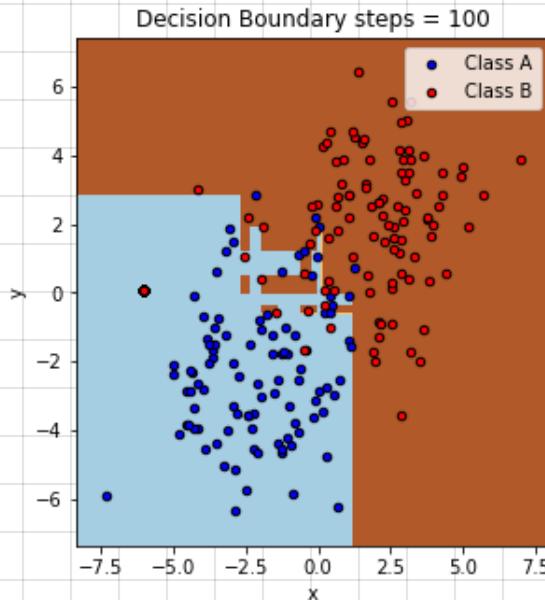
Boosting (Schapire and Freund 2012) is a greedy algorithm for fitting adaptive basis-function models of the form shown above where ϕ_m are generated by an algorithm called a weak learner or a base learner. One of the popular boosting algorithms is the *adaboost*.

For the data-set given [at this link](#), determine the decision boundary by using adaboost. Run the the algorithm for 250 steps. Show the decision boundary at steps 10, 100, 200 and 250. For additional details, refer to `demo_boosting` code in the PMTK package associated with Murphy [1].

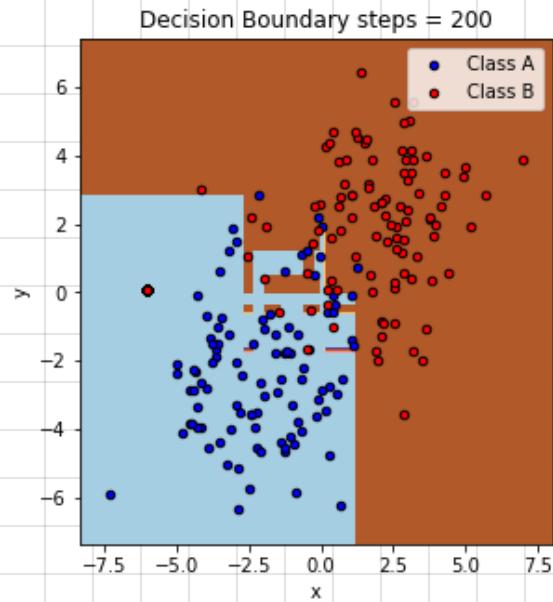
Steps = 10



Steps = 100



Steps = 200



Steps = 250

